



(12) **EUROPEAN PATENT APPLICATION**

(43) Date of publication:
23.12.1998 Bulletin 1998/52

(51) Int Cl.⁶: **G06F 17/27, G06F 17/30**

(21) Application number: **98304842.2**

(22) Date of filing: **19.06.1998**

(84) Designated Contracting States:
AT BE CH CY DE DK ES FI FR GB GR IE IT LI LU
MC NL PT SE
 Designated Extension States:
AL LT LV MK RO SI

(72) Inventor: **Julliard, Laurent**
38100 Grenoble (FR)

(74) Representative: **Reynolds, Julian David et al**
Rank Xerox Ltd
Patent Department
Parkway
Marlow Buckinghamshire SL7 1YL (GB)

(30) Priority: **20.06.1997 GB 9713019**

(71) Applicant: **XEROX CORPORATION**
Rochester, New York 14644 (US)

(54) **Linguistic search system**

(57) A method of searching for information in a text database, comprising: receiving (s1) at least one user input, the user input(s) defining a natural language expression, converting (s2, s3) the natural language expression to a tagged form (50, 51) including part-of-speech tags, applying (s4) to the tagged form (51) one or more grammar rules of the language of the natural language expression (49), to derive a regular expression (52), and analysing (s5) the text database to determine whether there is a match between said regular expression (52) and a portion of said text database. An apparatus for carrying out this techniques is also disclosed. Users may find portions of a text which match multiword expressions given by the user. Matches include possible variations that are relevant with the initial criteria from a linguistic point of view including simple inflections like plural/singular, masculine/feminine or conjugated verbs and even more complex variations like the insertion of additional adjectives, adverbs, etc. in between the words specified by the user.

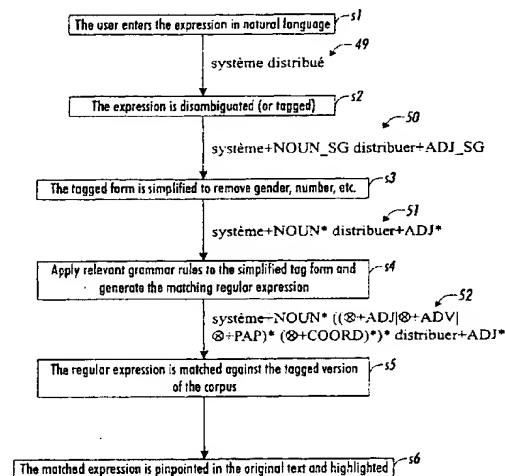


FIG. 2

Description

The present invention relates to data processing, and more particularly to techniques for searching for information in a text database or corpus.

Most of the techniques in use to retrieve a piece of information in a text corpus are based on substring search (also known as full-text search). Because this basic string search mechanism is weak when the user wants to catch more than a simple sequence of characters, various techniques have been developed by data providers to enhance the substring matching. Examples include wildcards, regular expressions, Boolean operators, proximity factor (e.g. words must be in the same sentence or no more than N words between two words) and stemming.

Existing techniques often try to achieve the similar goals: to allow the user to better express the variability of the natural language in which the string expression is to be searched in order not to miss any place where this expression appears.

However, known techniques suffer from several drawbacks: the end user has to learn the query language proposed by the search engine; no two search engines have the same query language; if the user doesn't think of all the possible variations of the searched expression, he can miss some relevant documents; and/or on the other hand, if the search expression is too "loose", many irrelevant documents will be retrieved, generating noise.

The present invention provides a method of searching for information in a text database, comprising: (a) receiving at least one user input, the user input(s) defining a natural language expression including one or more words, (b) converting the natural language expression to a tagged form of the expression, the tagged form including said one or more words and, for the or each word, a part-of-speech tag associated therewith, (c) applying to the tagged form one or more grammar rules of the language of the natural language expression, to derive a regular expression, and (d) analysing the text database to determine whether there is a match between said regular expression and a portion of said text database.

Preferably, step (b) comprises the step of: tagging the natural language expression by, for the or each word in said natural language expression, (b1) converting the word to its root form, and (b2) applying a part-of-speech tag to the word, thereby generating a complex tagged form.

Preferably, the part-of speech tag includes a syntactic category marker and a morphological feature marker, and wherein step (b) further comprises the step of: (b3) simplifying said complex tagged form by removing the or each morphological feature marker, to generate said tagged form.

Preferably, the method further includes the step of (e) determining the location of said text database of a

match with said regular expression.

The invention further provides a programmable data processing apparatus when suitably programmed for carrying out the method of any of the appended claims, or according to any of the particular embodiments described herein, the processor being coupled to the memory and user interface, and being operable in conjunction therewith for executing instructions corresponding to the steps of said method(s).

The linguistic search techniques according to the present invention overcome at least some of the above-mentioned problems. They rely both on the linguistic tools (such as atokeniser, morphological analyser and disambiguator) and the generation of complex regular expressions to match against the text database.

This mechanism has the advantages over a basic full text search engine that the end user doesn't need to learn an esoteric query language. He just has to type the multiword expression he is looking for in natural language.

A further advantage is that the retrieved documents will be much more relevant to the query from a linguistic point of view (although it doesn't ensure that all relevant documents will be retrieved from the point of view of the meaning).

A further advantage is that many variations will be captured by the linguistic processing. As a consequence, even a user who is not familiar with the language in which the searched documents are written doesn't have to know about the linguistic variation that might occur.

The linguistic search techniques according to the invention provide a new way to search for information in a text database. They enable users to find portions of a text which match multiword expressions given by the user. Matches include possible variations that are relevant with the initial criteria from a linguistic point of view including simple inflections like plural/singular, masculine/feminine or conjugated verbs and even more complex variations like the insertion of additional adjectives, adverbs, etc. in between the words specified by the user. This technique can complement conventional full text search engines by reducing the number of retrieved documents that are inconsistent with the query.

Embodiments of the invention will now be described, by way of example, with reference to the accompanying drawings, in which:

Figure 1 is a schematic block diagram of the computer which may be used to implement the techniques according to an embodiment of the present invention; and

Figure 2 is a schematic flow diagram of the steps in carrying out a linguistic search according to an embodiment of the present invention.

It will be appreciated that the present invention may be implemented using conventional computer technol-

ogy. The invention has been implemented in Perl & C++ on a Sun workstation running SunOS. It will be appreciated that the invention may be implemented using a PC running Windows™, a Mac running MacOS, or a minicomputer running UNIX, which are well known in the art. For example, the PC hardware configuration is discussed in detail in *The Art of Electronics*, 2nd Edn, Ch. 10, P. Horowitz and W. Hill, Cambridge University Press, 1989, and is illustrated in Fig. 1. Stated briefly, the system comprises, connected to common bus 30, a central processing unit 32, memory devices including random access memory (RAM) 34, read only memory (ROM) 36 and disk, tape or CD-ROM drives 38, keyboard 12 (not shown), mouse 14 (not shown), printing, plotting or scanning devices 40, and A/D, D/A devices 42 and digital input/output devices 44 providing interfacing to external devices 46 such as the rest of a LAN (not shown).

Figure 2 is a schematic flow diagram of the steps performed in carrying out a linguistic search according to an embodiment of the present invention.

It will be apparent to persons skilled in the art that where references are made herein to steps, operations or manipulations involving characters, words, passages of text, etc., these are implemented, where appropriate, by means of software controlled processor operations upon machine readable (e.g. ASCII code) representations of such characters, words and text.

For the sake of illustrating the techniques according to the invention, the case is considered where the French expression "système distribué" (equivalent to "distributed system" in English) is to be searched for in a French corpus by the user.

Initially (step s1), the user specifies the multiword expression he is looking for, for example, using the type of graphical user interface which is well known in the art. There is no need to pay attention to the formulation of this expression: nouns and/or adjectives can be plural or singular, verbs can be conjugated, etc.

Next, at step s2, the expression is then sent to the tagger (or disambiguator), such as are available from Xerox Corp. Taggers are discussed in more detail in McEnery T. and Wilson A., *Corpus Linguistics*, Ch. 5, section 3 and Appendix B. The tagger (or disambiguator) does two things—

- (1) reduce each word to its root form (e.g. *distribué* becomes *distribuer* - infinitive form of the verb), and
- (2) determine the part-of-speech of each word (e.g. *système* is a singular noun - NOUN_SG- and *distribué* is a singular adjective -ADJ_SG-). NOUN_SG and ADJ_SG are called tags. Each tag is made of two parts: the syntactic category (or part-of-speech like NOUN, ADJ, VERB, etc.) and the morphological feature (like SG, PL, etc.) which reflects the inflection of the word.

Once the tagged form 50 has been obtained, it is then simplified, at step s3: because it is desired that the

linguistic search process retrieves all possible inflections of a word each tag is first reduced to its syntactic category. The gender, number or person of a word is useless for the linguistic search, and is removed. Preferably, this comprises replacing each of "SG", "PL", etc. with a neutral symbol (*) so as to encompass all possibilities of morphological feature.

The process continues at step s4, in which the simplified tagged form 51 is operated on. Given the grammar of a language it is possible to determine what kind of variations a multiword expression can undergo without changing its initial meaning. The following discussion presents some of the rules that have been used for French to generate variations around nominal phrases:

- (1) In between a noun and an adjective one can insert adjectives, adverbs, or past participles possibly connected by a co-ordinating conjunctions like *et* (*and*), *ou* (*or*), etc. Figure 2 illustrates the application of this rule to the expression *systèmes distribués* and shows a simplified version of the resulting regular expression (the symbol ® represents the word preceding the tag). As an example, there follow some linguistic variations caught by this regular expression:

⇒ *systèmes distribués* (*distributed systems* - plural form)
 ⇒ *systèmes relationnels distribués* (*distributed relational systems* - inserted adjective)
 ⇒ *système redondant et totalement distribué* (*fully redundant and distributed system* - inserted adjective and adverb joined by a co-ordinating conjunction)

- (2) In between a noun and a preposition or a preposition and a noun, additional adjectives can be inserted.

- (3) In between 2 nouns, additional adjectives can be inserted.

The rules listed above apply to French noun phrases. They can be extended to any other kind of phrases, including those containing verbs, and also to any other language.

It is to be noted that these rules can be almost as complex as desired if it is thought that there is a good chance for the selected portion of text to be still relevant with your initial query. For instance one could allow the insertion of a new noun phrase in between the noun and the adjective like in "*système à tolérance de panne distribué*" (*distributed fault tolerant system*) or even more complex is the insertion of a relative clause like in "*un système qui, par nature, est totalement distribué*" (*a system which, by essence, is fully distributed*).

The grammar rules expressed in step s4 are coded in a regular expression and matched against the simplified tagged form 51 of the user query. If one of those

rules matches, then the simplified tagged form 51 of the user query transforms into a complex regular expression representing the grammar variations.

Each rule is applied in sequence and only once in order to avoid the recursive application of a grammar rule to itself or to others.

The matching regular expression 52 is then processed further at step s5. Once the final regular expression 52 has been generated it is matched against the tagged version of the corpus. With respect to this step, it is important to note the following.

(1) As stated above the matching process has to be made on a tagged version of the text corpus. This may be done using a tagger, such as that available from Xerox Corp., as mentioned above. The tagging phase can be made either on the fly, if the text varies frequently, or once for all if it is stable.

(2) If the corpus is large, a simple sequential search on the tagged text will take too much time. To speed up this phase a full text indexing engine can be used. But instead of indexing the original text as most full-text search engines do, the indexing mechanism is applied to the tagged version of the text corpus.

(3) Most existing full-text indexing engines cannot handle search queries expressed with complex regular expressions. As a consequence the expression generated by the linguistic search system according to the present invention cannot be given as is to the search engine. In fact, a preliminary search is made on the individual words of the simplified tagged expression (see step s2). Depending on how sophisticated the indexing engine is, it can provide the user with very basic information like the name of the files in which those words have been found (like the glimpse search engine does) or much more accurate information like the position of the sentence in which those words were found (like the Xerox PARC Text Database (TDB) does). This preliminary step reduces the scope of relevant (portion of) documents and decreases the time required by the regular expressions matching process.

(4) The current implementation of an embodiment of the linguistic search system according to the invention is based on the regular expression conventions of Perl (or any flavour of awk). It will be appreciated by persons skilled in the art that it could be easily transposed to the regular expression formalism used by the Finite State Transducers developed by Xerox Corp. (see EP-A-583,083). The matching mechanism is based on the regular expressions of Perl rather than the Finite State transducers developed by Xerox because Perl (and awk) tells the user not only what portion of the text matched but also where it is located in the corpus. This information is particularly noticeable in order to highlight the places where a match occurred. This feature has two

advantages:

(1) avoid leafing through long documents to find places where matches occurred (See step s6 discussed below); and

(2) show the whole matching multiword expression which can be quite different from the one typed by the user if the linguistic variations allowed by the grammar rules are complex.

Step s6 is performed after the regular expression has been matched against the tagged version of the corpus. As mentioned above, the Perl (or awk) regular expressions mechanism can tell the user what string matches but also where this string is located in the text. However because according to the invention the regular expression matching is done on the tagged version of the corpus, the positioning information is not suitable for the original text. As a consequence, if it is desired to highlight the matches a way must be provided to go from the offset in the tagged text into the actual offset in the original text. Currently, this is made via a simple offset table built during the corpus tagging.

It will be appreciated that numerous modifications may be made in implementing the techniques according to the invention.

The linguistic search could be applied to WEB search engines. Although their query languages tend to be more and more sophisticated it's not yet close to a linguistic search.

The process explained above assumes that the corpus to be searched is first disambiguated (or tagged). However, it will be appreciated that it would be possible to use the techniques according to the invention as a front-end to the WEB search engine, for example. Here, the requirement is to generate all the possible forms of a word and search for all of them with a conventional search engine (or at least the substring which is common to all the derived form of a word). Then the selected documents need to be retrieved for further processing (tagging) before the linguistic search can be applied.

4) References

1. LOCOLEX: Translation Rolls off Your Tongue. Daniel Bauer, Frédérique Segond and Annie Zaenen, RXRC, Grenoble, FRANCE, in the *Proceedings of the conference of the Association for Computers and the Humanities and the Association for Literary and Linguistic Computing*, (ACH-ALLC'95) Santa Barbara, USA, July 1995.
2. SEXTANT: Extracting Semantics from Raw Text. Gregory Grefenstette, RXRC, Grenoble, FRANCE, in *Integrated Computer-aided Engineering* July 1993.
3. Constructing Lexical Transducers. Lauri Karttunen, RXRC Grenoble FRANCE, in *COLING'94 Proceedings*.

4. Creating a tagset, lexicon and guesser for a French tagger. Jean-Pierre Chanod and Pasi Tapanainen, in *Proceedings of ACL-SIGDAT*, 1995.
5. Retrieving terms and their Variants in a Lexicalized Unification-Based Framework. Christian JACQUEMIN and Jean ROYAUTE, in *Proceedings of ACM-SIG Information Retrieval*, July 1994.
6. Automatic Search Term Variant Generation. K. Sparc Jones, Computer Laboratory, University of Cambridge, UK in *Journal of Documentation*, Vol. 40, No. 1, March 1984, pp. 50-66.
7. Natural Language Processing: the PLNLP Approach. Karen Jensen, George E. Heirdon, Stephen D. Richardson. Microsoft Corporation. KLUWER ACADEMIC PUBLISHERS.
8. Information Retrieval and Virtual Libraries: the Callimaque model. Monica Beltrametti, Laurent Julliard, Françoise Renzetti. *Proceedings of CAIS'95*.

4. The method of claim 1, 2 or 3, further including the step of (e) determining the location of said text database of a match with said regular expression (52).
5. A programmable data processing apparatus when suitably programmed for carrying out the method of any of the preceding claims, the apparatus including a processor, memory, and a user interface, the processor being coupled to the memory and user interface, and being operable in conjunction therewith for executing instructions corresponding to the steps of said method(s).

Claims

1. A method of searching for information in a text database, comprising:
 - (a) receiving at least one user input, the user input(s) defining a natural language expression (49) including one or more words,
 - (b) converting the natural language expression to a tagged form (50, 51) of the expression, the tagged form including said one or more words and, for the or each word, a part-of-speech tag associated therewith,
 - (c) applying to the tagged form (51) one or more grammar rules of the language of the natural language expression (49), to derive a regular expression (52), and
 - (d) analysing the text database to determine whether there is a match between said regular expression (52) and a portion of said text database.
2. The method of claim 1, wherein step (b) comprises the step of:
 - tagging the natural language expression by, for the or each word in said natural language expression, (b1) converting the word to its root form, and (b2) applying a part-of-speech tag to the word, thereby generating a complex tagged form (50).
3. The method of claim 2, wherein the part-of speech tag includes a syntactic category marker and a morphological feature marker, and wherein step (b) further comprises the step of: (b3) simplifying said complex tagged form (50) by removing the or each morphological feature marker, to generate said tagged form (51).

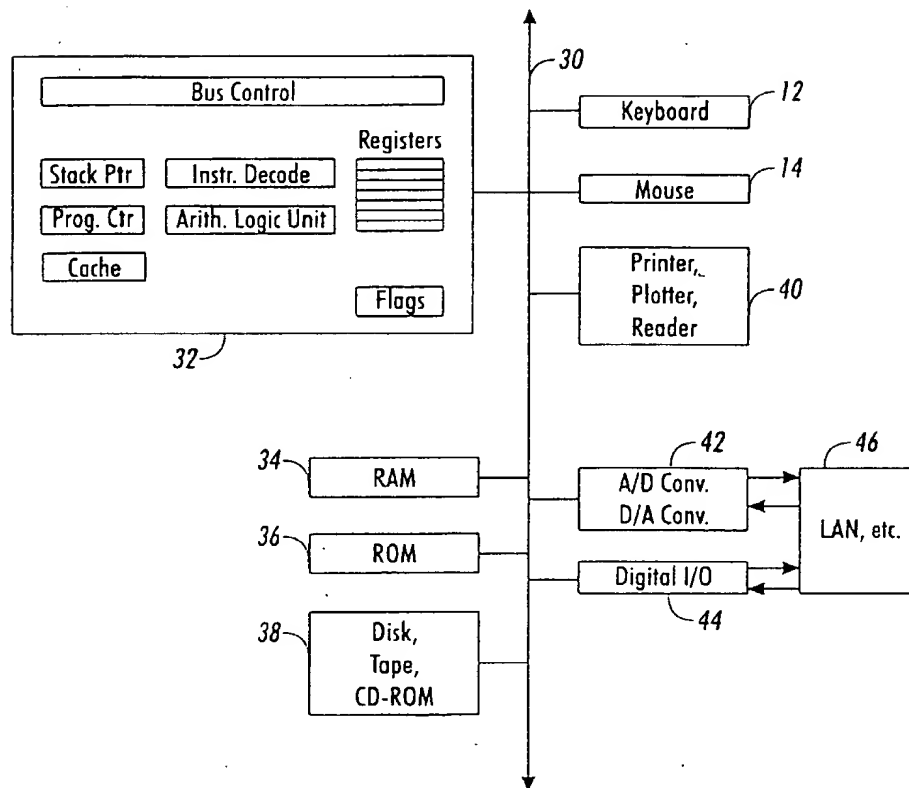


FIG. 1

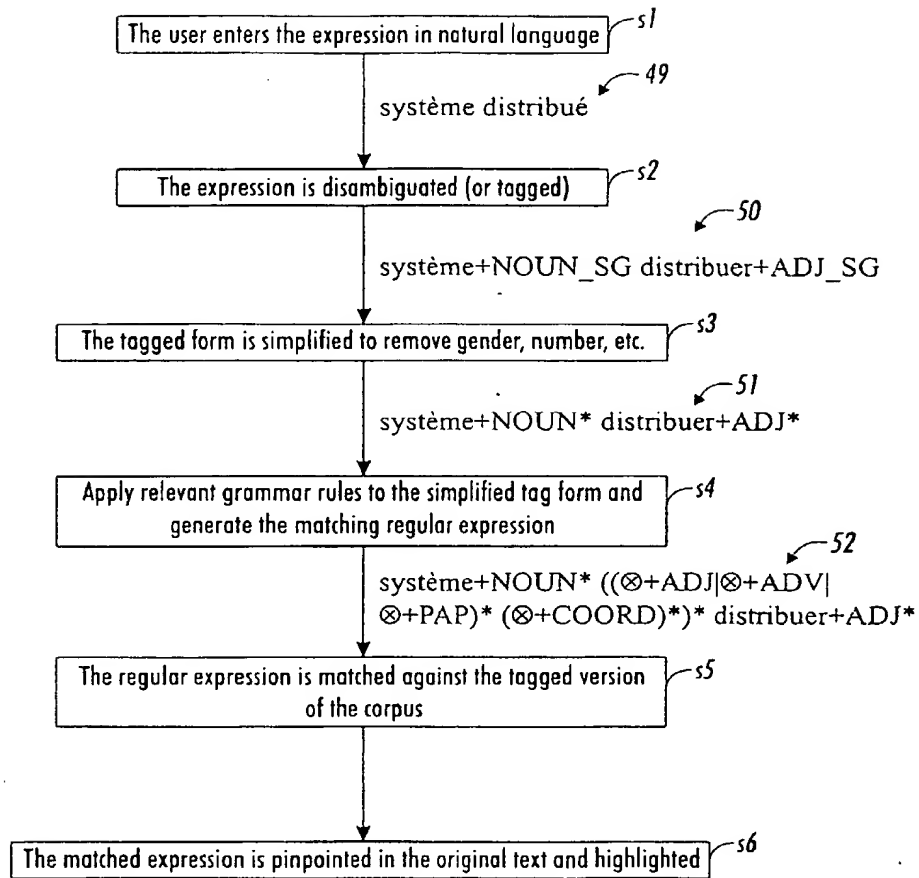


FIG. 2



European Patent
Office

EUROPEAN SEARCH REPORT

Application Number
EP 98 30 4842

DOCUMENTS CONSIDERED TO BE RELEVANT			
Category	Citation of document with indication, where appropriate, of relevant passages	Relevant to claim	CLASSIFICATION OF THE APPLICATION (Int.Cl.6)
X	ANTONIADIS G ET AL: "A FRENCH TEXT RECOGNITION MODEL FOR INFORMATION RETRIEVAL SYSTEM" PROCEEDINGS OF THE INTERNATIONAL CONFERENCE ON RESEARCH AND DEVELOPMENT IN INFORMATION RETRIEVAL. (SIGIR), GRENOBLE. JUNE 13 - 15, 1988. no. CONF. 11, 13 June 1988, pages 67-84, XP000246106 CHIARAMELLA Y * page 67, line 12 - line 32 * * page 74, line 20 - page 83, line 14 *	1-5	G06F17/27 G06F17/30
A	EP 0 597 630 A (CONQUEST SOFTWARE INC) 18 May 1994 * page 7, line 46 - page 10, line 43 *	1-5	
A	US 5 418 716 A (SUEMATSU HIROSHI) 23 May 1995 * column 5, line 8 - column 6, line 56 *	1-5	
A	PATENT ABSTRACTS OF JAPAN vol. 015, no. 180 (P-1199), 9 May 1991 & JP 03 040067 A (NEC CORP). 20 February 1991 * abstract *	1-5	TECHNICAL FIELDS SEARCHED (Int.Cl.6) G06F
The present search report has been drawn up for all claims			
Place of search BERLIN		Date of completion of the search 6 October 1998	Examiner Deane, E
<p>CATEGORY OF CITED DOCUMENTS</p> <p>X : particularly relevant if taken alone Y : particularly relevant if combined with another document of the same category A : technological background O : non-written disclosure P : intermediate document</p> <p>I : theory or principle underlying the invention E : earlier patent document, but published on, or after the filing date D : document cited in the application L : document cited for other reasons & : member of the same patent family, corresponding document</p>			

EPO FORM 1503 03/82 (F04C01)